# Data Publication Implementation and Metrics

Deborah Smith

**2017 GHRC User Working Group Meeting**
Sept 26-27, 2017

# What is DAPPeR?



DAPPeR is an online publication tool based on ORNL SAuS

- Put into operation in May 2017

DAPPeR web page contains description of publication process and a direct link to the data archival interest form

- Made public in August 2017

The secure portal interface has access for both data providers and GHRC personnel

DAPPeR URL: https://ghrc.nsstc.nasa.gov/data-publication/

# DAPPeR: Dataset Submission on Web

- Description of publication process is public with needed links for more information and form access

- Feedback link is available for questions/comments

## HOW TO PUBLISH YOUR DATA AT GHRC

The entire workflow of dataset publication at GHRC from accepting it for archival to publishing it to the end user is managed by the GHRC Data Publiction Portal and associated Data Publication Portal used by the GHRC DMG. The portal provides a centralized system for Data Producers to provide dataset details, for GHRC DMG and GHRC DAAC Managers to push the dataset through the stages of publication, and for the GHRC DMG to interact with the data producers about the datasets. The portal automates many of the steps required for dataset publication to reduce redundancy and repetition of effort, and to speed the overall process. Data providers can monitor their submissions as DMG works toward data publication.

When you submit data to GHRC, you agree to work with the team in a timely manner to assist with the process of data publication.

**You may contact GHRC at any time during the process using this data product submission portal.**

This detailed workflow chart shows how the GHRC ingests and publishes your data through various stages, what team members handle your data and what your responsibilities are to the process.

**Step 1:** Login with your Earthdata account. If you do not have an Earthdata login account, learn how to create one here: https://ghrc.nsstc.nasa.gov/home/access-data/earthdata-login-recipe

**Step 2:** Complete and submit the Archival Interest Form

**Step 3:** When you obtain the news that GHRC is interested in publishing your data, complete and submit the PI Questionnaire Form. If you have any questions, you can use the portal to contact the DMG during the process.
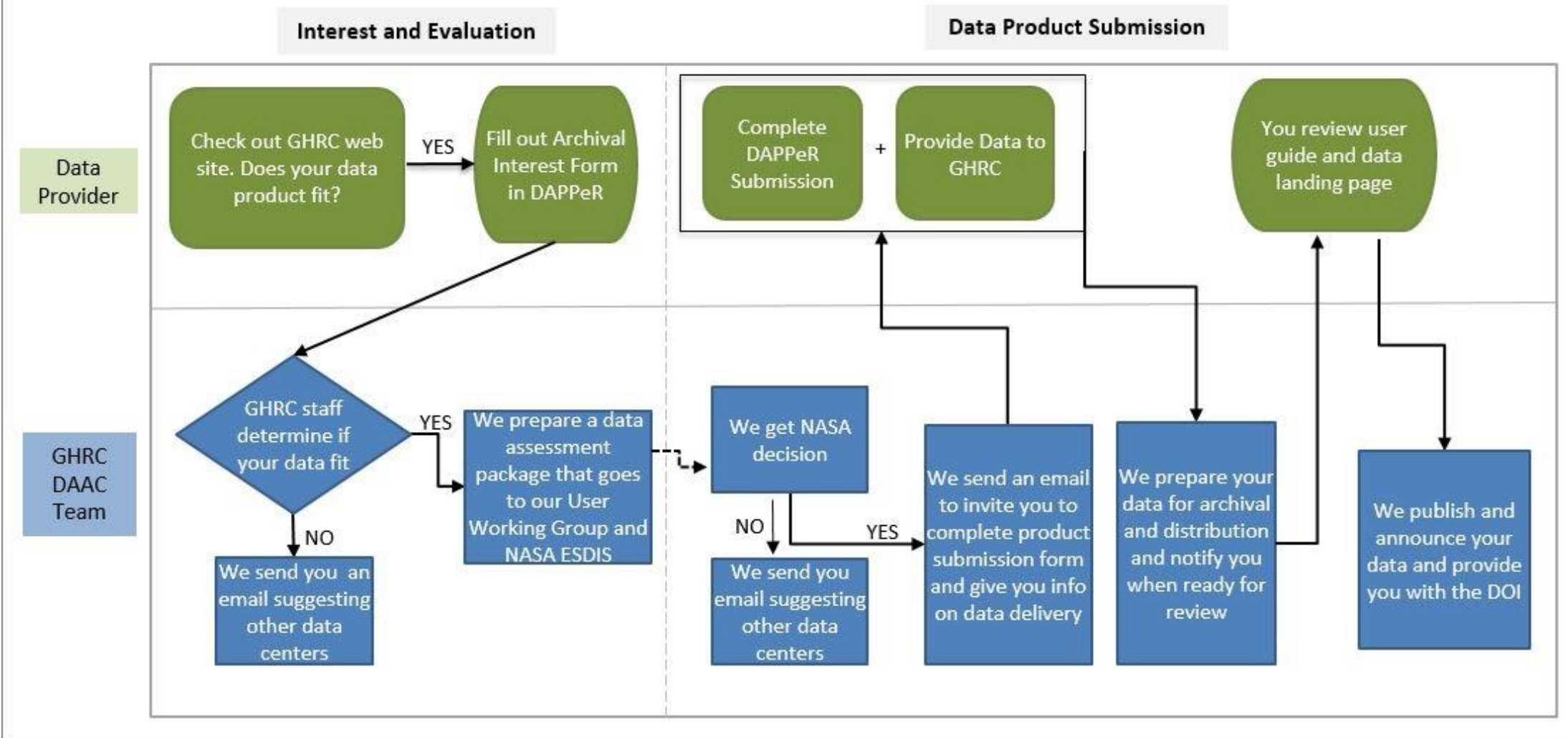
**Step 4:** Submit your data following the instructions provided to you by the DMG.

**Step 5:** The DMG will next catalog and process your data, write documentation to support your data, and create browse images of your data if necessary. We will contact you if there are any questions.

**Step 6:** You will be contacted by the DMG for feedback on the documentation produced. Let the team know as soon as possible that all is well.

**Step 7:** The DOI and dataset landing pages are generated and activated as your dataset is published at the GHRC DAAC. At this point you can refer others to the public data and provide the DOI to journals as needed.

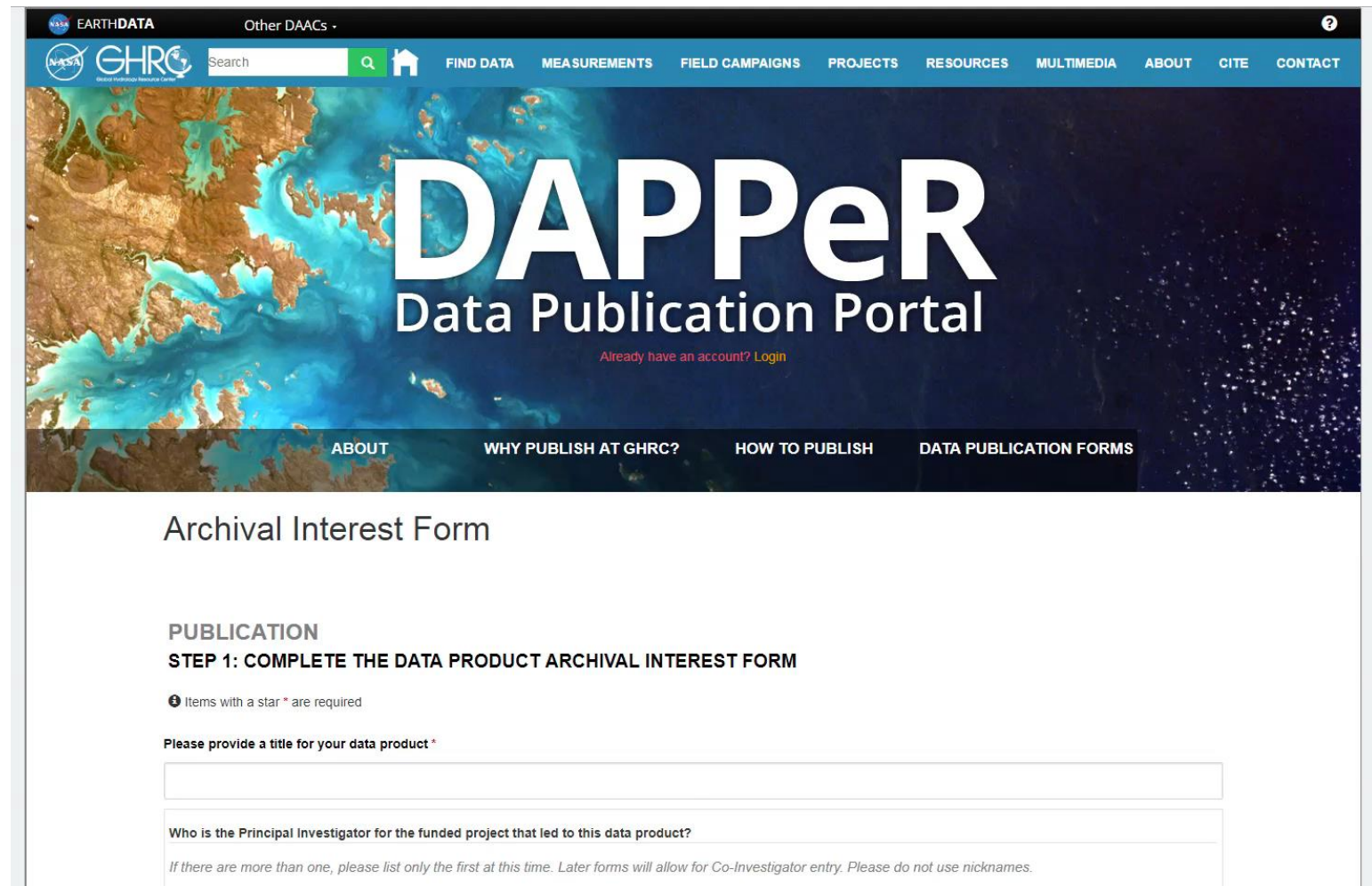How do you publish data at GHRC?

# Data Publication Forms

- **Archival Interest Form**
  - Used to collect information about and relevancy of a proposed dataset to GHRC
  - https://ghrc.nsstc.nasa.gov/data-publication/node/add/archival_interest_form

- **Data Product Questionnaire**
  - Secure access for data providers to submit detailed dataset information once the dataset has cleared the review process
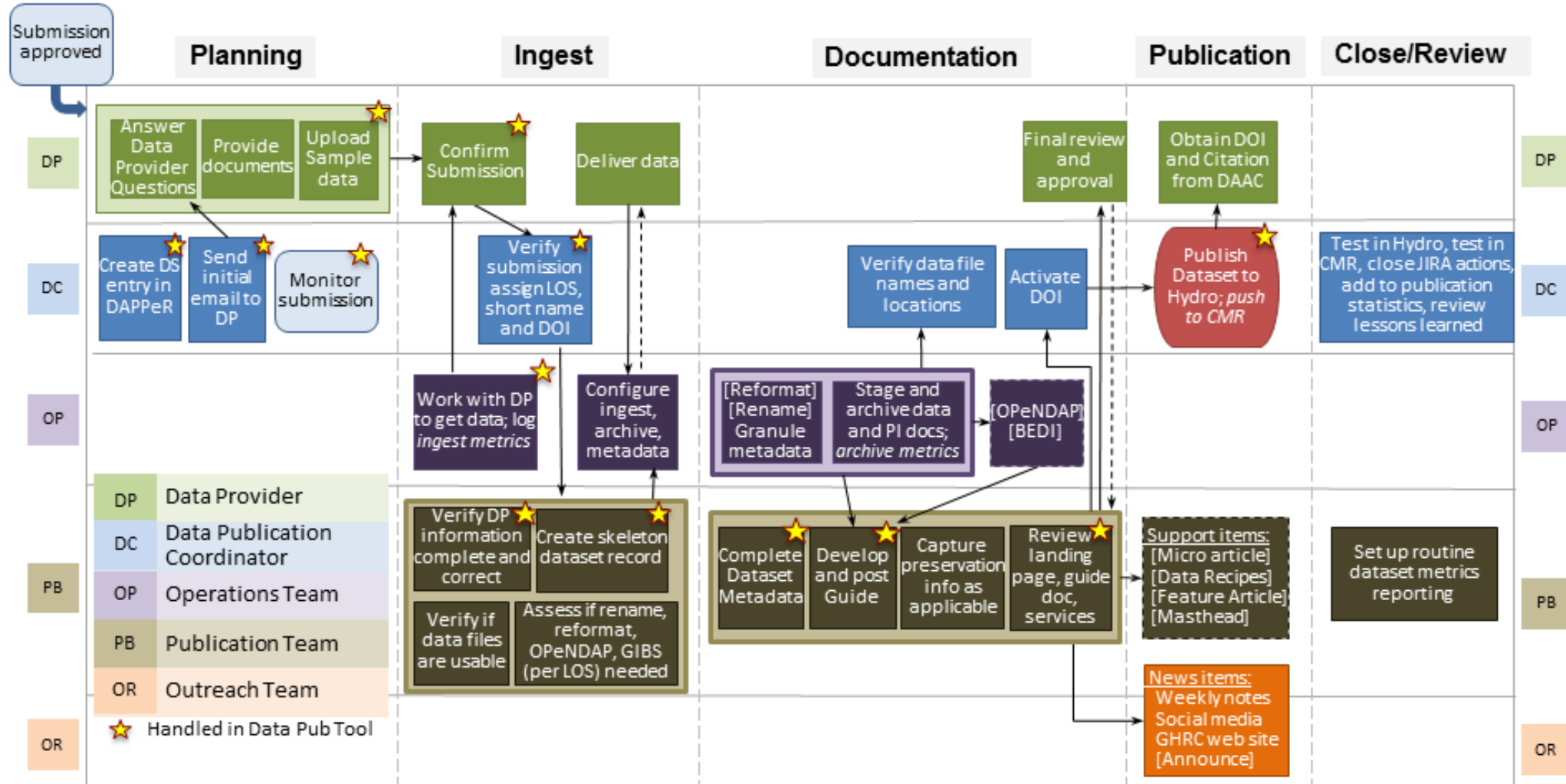  - Data producers are required to complete most questions in order for dataset to be accepted

# GHRC Data Publication Changes

**The GHRC data publication process underwent great change in 2017**

- A new database and publication tool DAPPeR was installed, tested, and implemented  during the beginning of 2017

- The publication team had to be trained in DAPPeR, including clear understanding of metadata definitions so that consistent information would be entered into the database by everyone

- During Jan - Apr, the publication team pre-prepared datasets

- Operational use of DAPPeR started in May 2017

- In July, Marilyn retired, resulting in a change of leadership

- In 4.5 months, we can already see an increase in publication rate from the new tools and processes
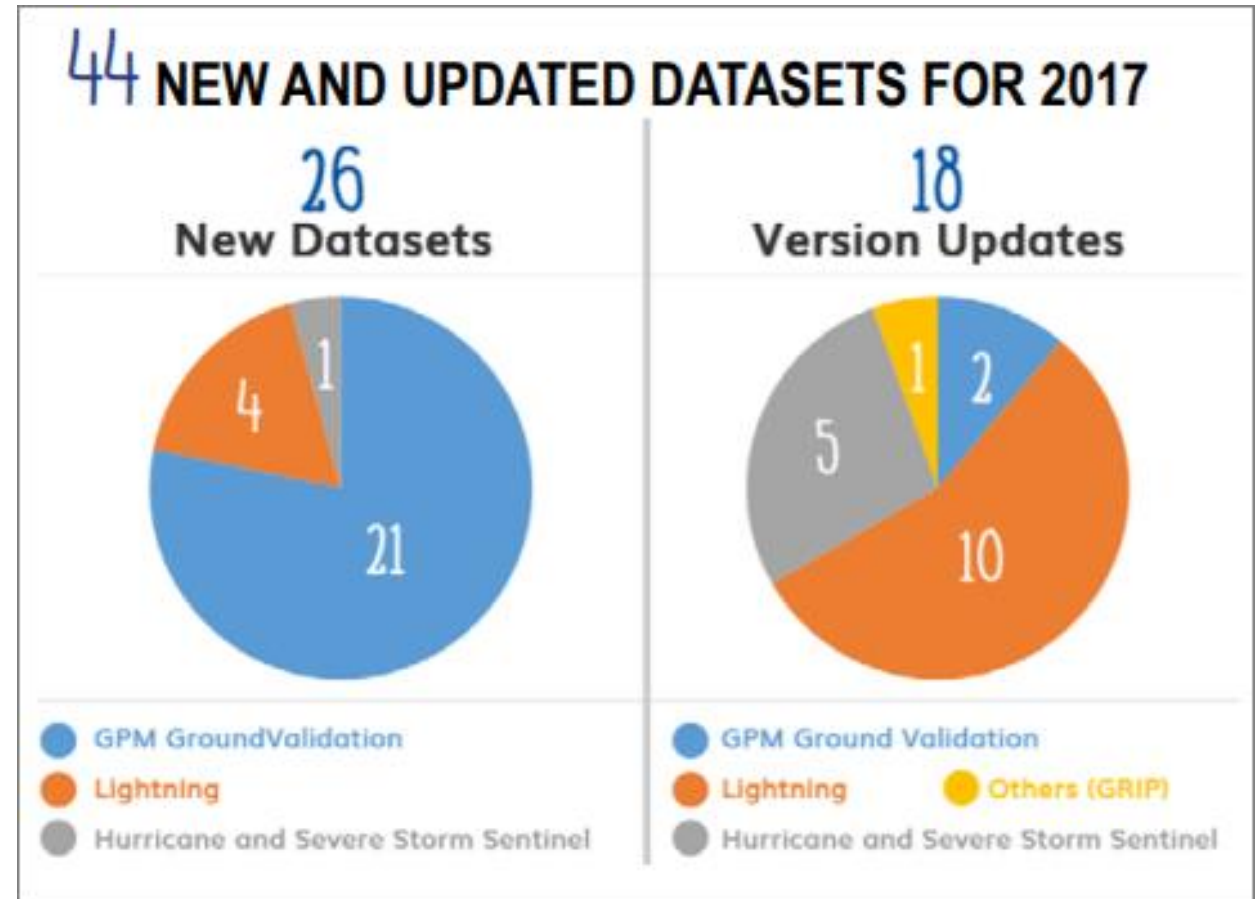
GHRC Data Publication Swimlanes

Progress has been made on publishing datasets already delivered to GHRC

Focus has been on first publishing the highest level of service datasets

- ISS LIS NRT and Non-QC datasets published

- OLYMPEX GPM Ground Validation
  - 50% of OLYMPEX campaign published

- All of HS3 core datasets are now published



44 NEW AND UPDATED DATASETS FOR 2017

26 New Datasets

18 Version Updates

New Datasets pie chart values: 21, 4, 1
- GPM GroundValidation
- Lightning
- Hurricane and Severe Storm Sentinel

Version Updates pie chart values: 2, 10, 5, 1
- GPM Ground Validation
- Lightning
- Others (GRIP)
- Hurricane and Severe Storm Sentinel

Comparison of dataset publication monthly totals for Jan - Sep 2016 and 2017



GHRC Monthly Data Publication Totals

Jan - Sep Totals
▲ 2016    29
● 2017    41

# Monthly Data Publication Rates Already Improved

## Comparison of dataset publication cumulative totals per year, 2016 -2018

**GHRC Monthly Data Publication Totals**

Jan - Sep Totals
- ▲ 2016  29
- ● 2017  41
- ◆ 2018  62

Cumulative Total Datasets Published

2018
2017
2016

Jan - April | May | June | July | August | September

# Publication Statistics From DAPPeR

- By using DAPPeR we will be able to record effort and time required for each dataset publication

- We anticipate to have statistics such as the following for next year's report:
    - avg time from start of effort to publication
    - avg time from ingest of dataset to publication
    - avg time for PI questionnaire completion
    - avg time per each step of publication
    - analyses by year, quarter, month
    - summary info on email communications, such as how many datasets required extra contact for publication

# HyDRO 2.0

- At last UWG meeting, we demonstrated our new HyDRO 2.0
    - Tool used to locate and access data from GHRC archive
    - New features made it easier for users to narrow down searches and find what they want
- HyDRO 2.0 was made fully operational in 2017
    - A YouTube instructional Video is provided to instruct on how to use HyDRO 2.0
- We incorporated a 'Get Help' button
- One single issue remains - Users having trouble downloading multiple files of data
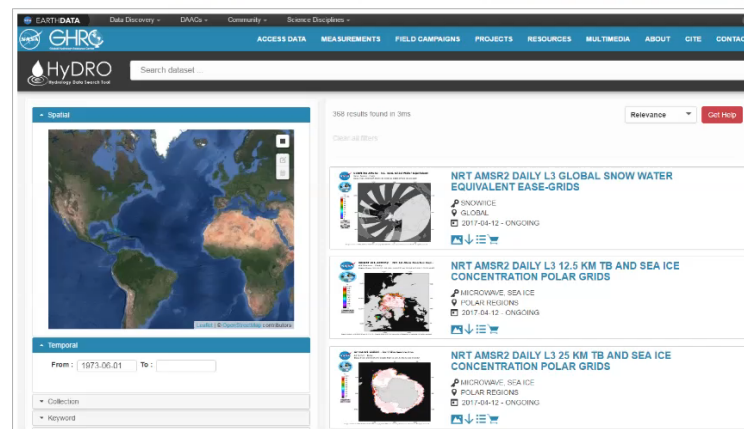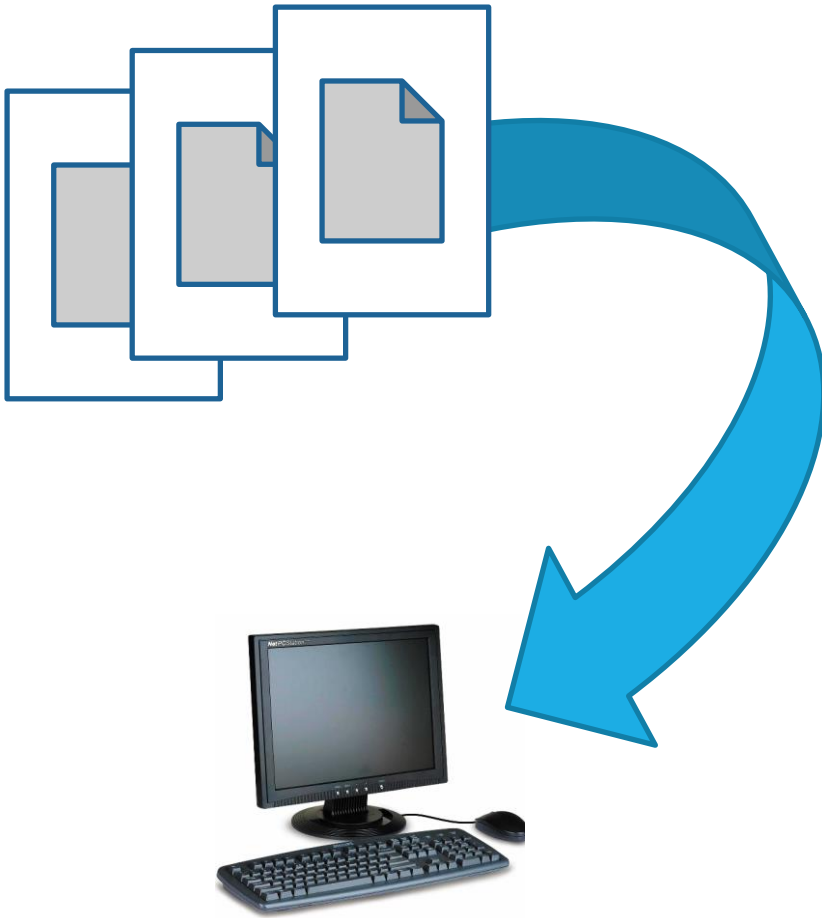
# Bulk Data Download Issues at GHRC

- Current process:  Utilize WGET with a provided URL file to bulk download data via HTTPS

- A text file of URLs is downloaded by the user thru HyDRO 2.0, then user runs  WGET calling the script
  - positives: works with Windows and Linux/Unix
  - negatives: requires user to install WGET

- WGET asks the user for Earthdata credentials before data download

- We are currently investigating other ways to bulk download data

- We have looked into Earthdata Drive, but had complications

- We have examined what other NASA DAAC's have for their own bulk data downloading tool.  Various approaches exist

- We are interested in finding a way to converge all of the diverse NASA DAAC bulk data download to create one standard bulk data download tool that can be used with Earthdata login across all DAACs

# Questions?

**2017 GHRC User Working Group Meeting**
Sept 26-27, 2017